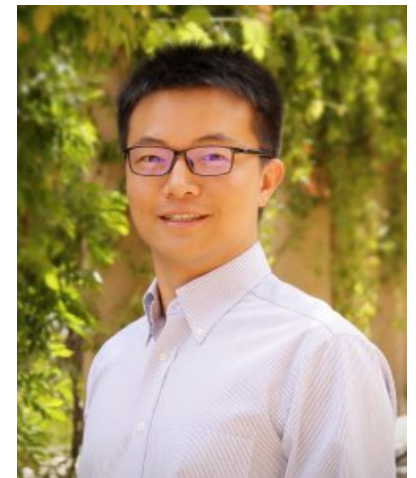


Compression in Language Modeling

Team: MIT-Han-Lab

Demi Guo^{*2}, Hanrui Wang^{*1}, Zhongxia Yan^{*1}, Phillip Isola¹ and Song Han¹



¹Massachusetts Institute of Technology

²Harvard University

^{*}Equal contribution, alphabetical order

Observations

- Desirable model properties
 - Predictive (< 35 PPL)
 - Fast to train
 - Efficient inference
- **Score = param/159M + compute/318M**
 - No penalty for non-parametric memory
 - Rely more on memory!

LSTM

- Sequential in time, slow to train
- Harder to learn long-term dependencies

This

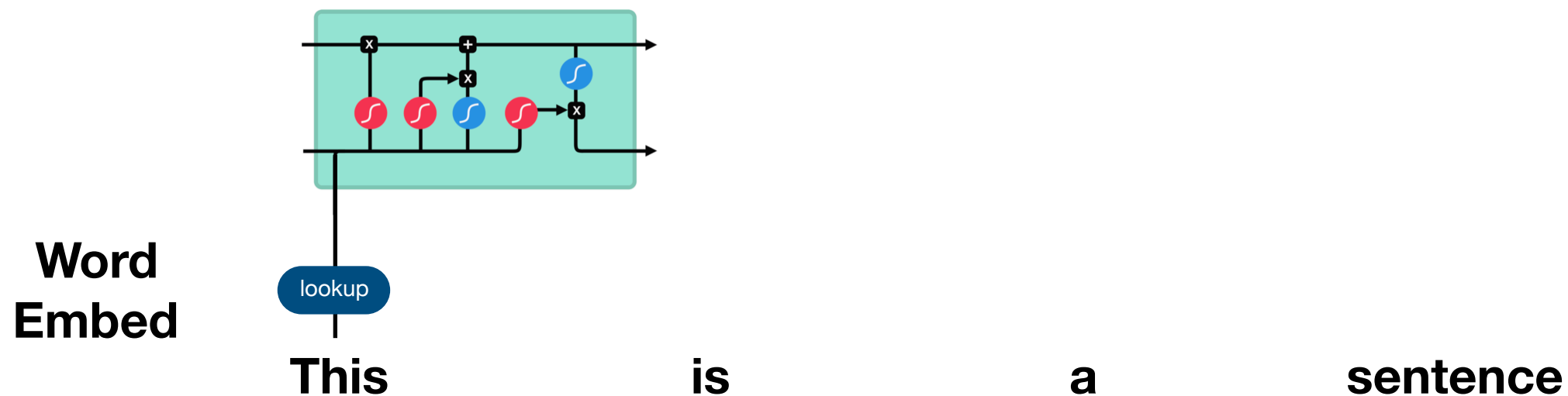
is

a

sentence

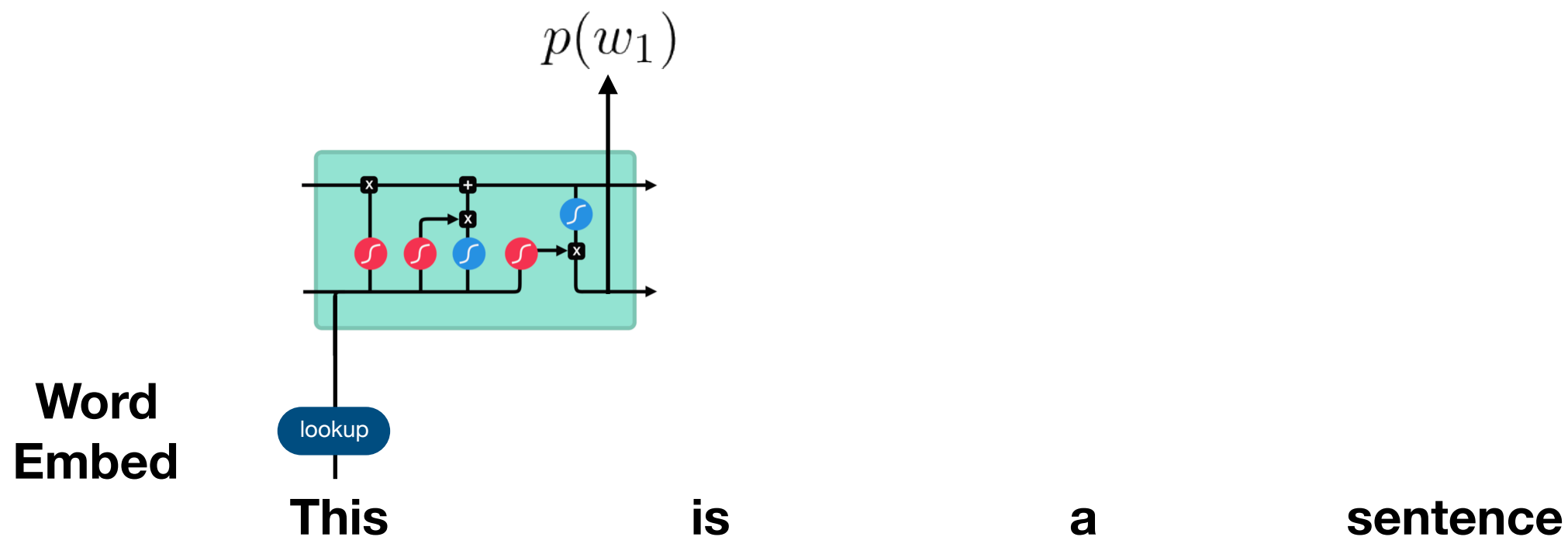
LSTM

- Sequential in time, slow to train
- Harder to learn long-term dependencies



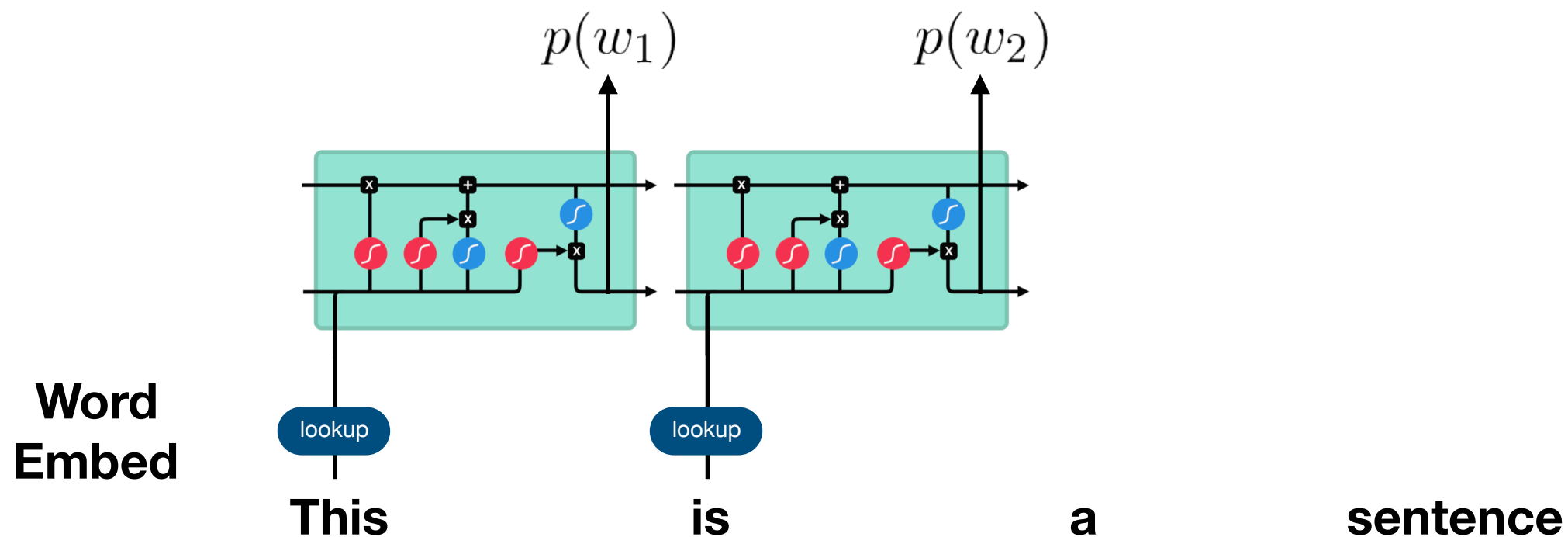
LSTM

- Sequential in time, slow to train
- Harder to learn long-term dependencies



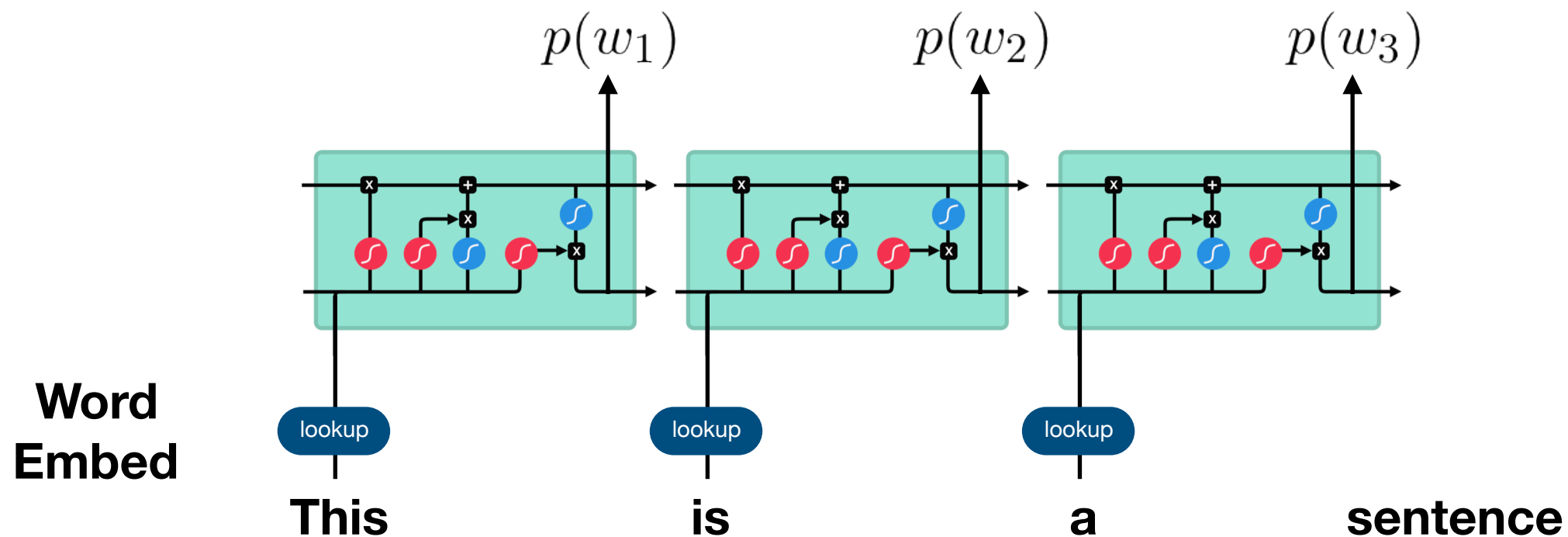
LSTM

- Sequential in time, slow to train
- Harder to learn long-term dependencies



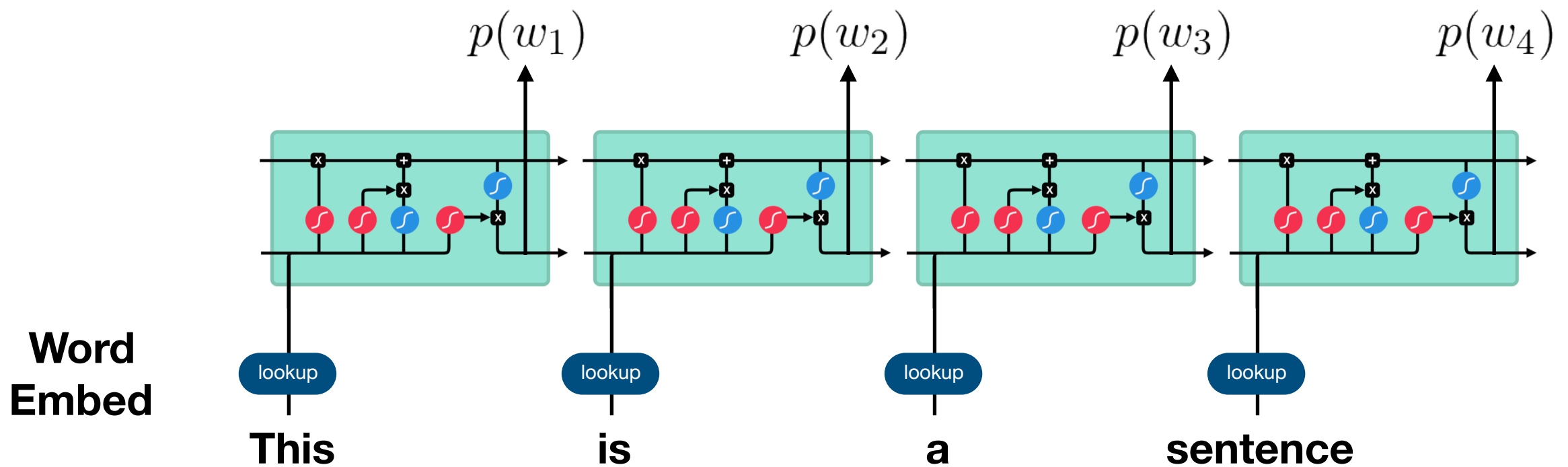
LSTM

- Sequential in time, slow to train
- Harder to learn long-term dependencies



LSTM

- Sequential in time, slow to train
- Harder to learn long-term dependencies



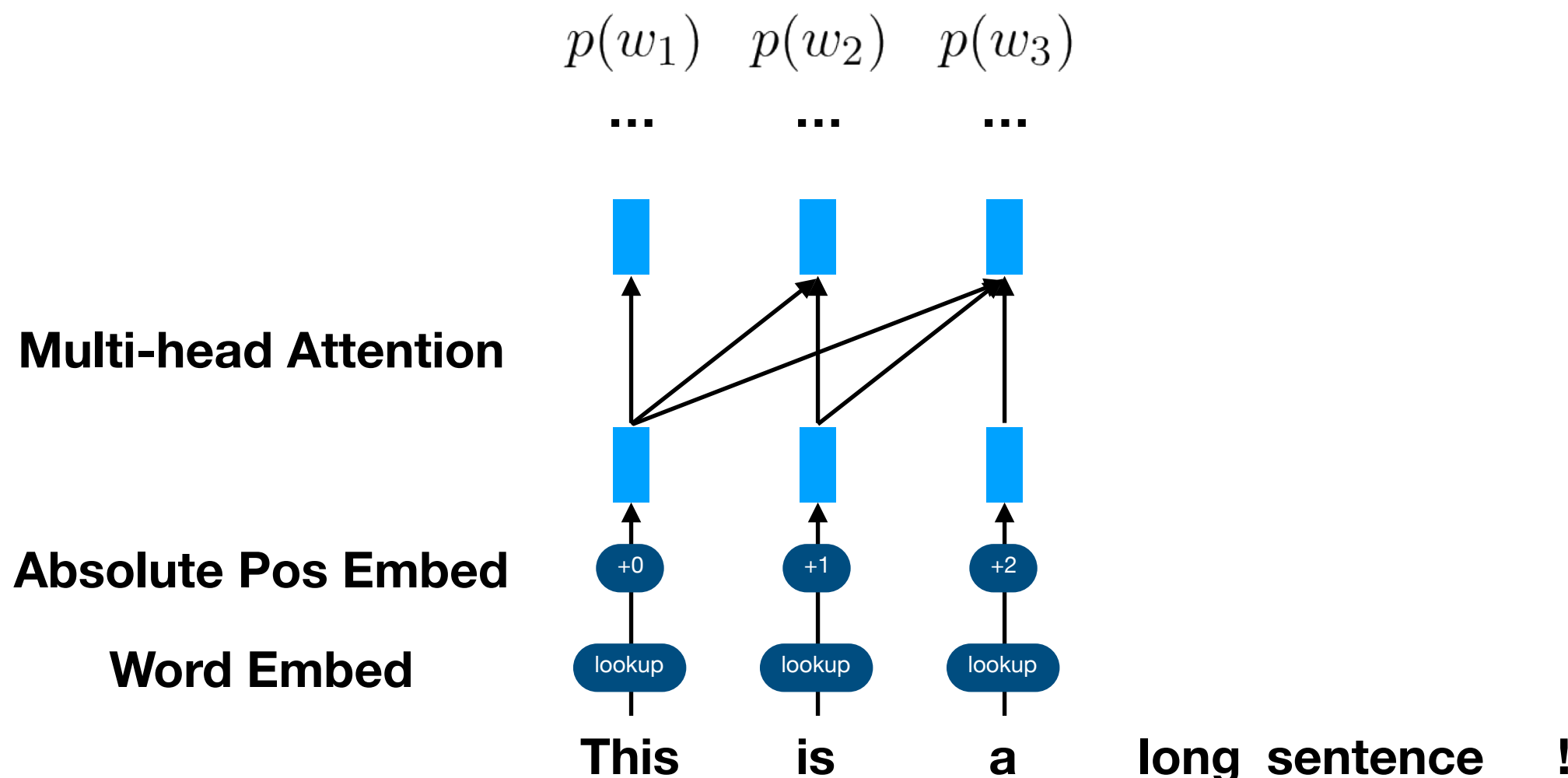
Transformer

- Transformer attention is parallel in time
- Absolute position embedding → recomputation during inference

This is a long sentence !

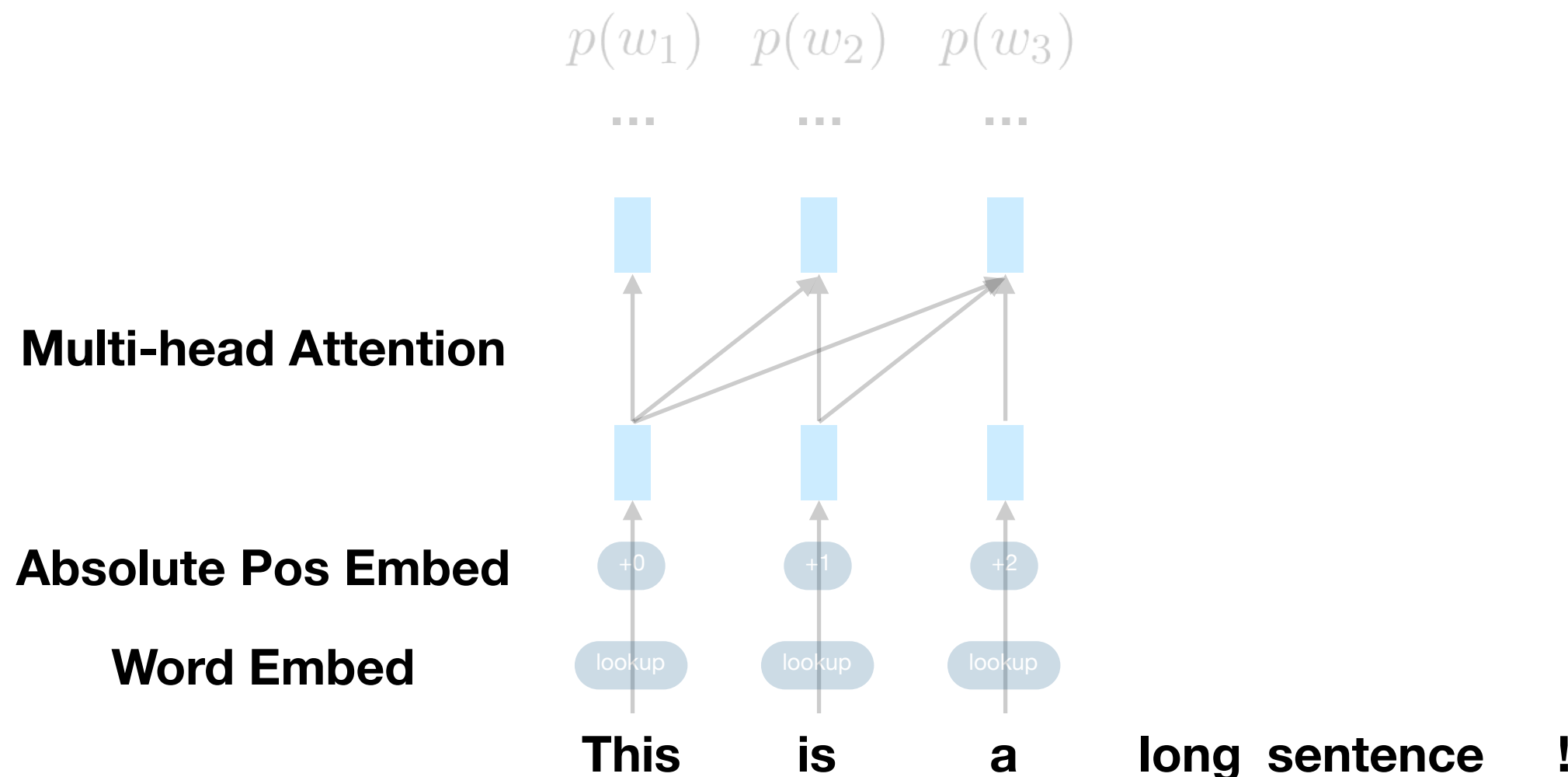
Transformer

- Transformer attention is parallel in time
- Absolute position embedding \rightarrow recomputation during inference



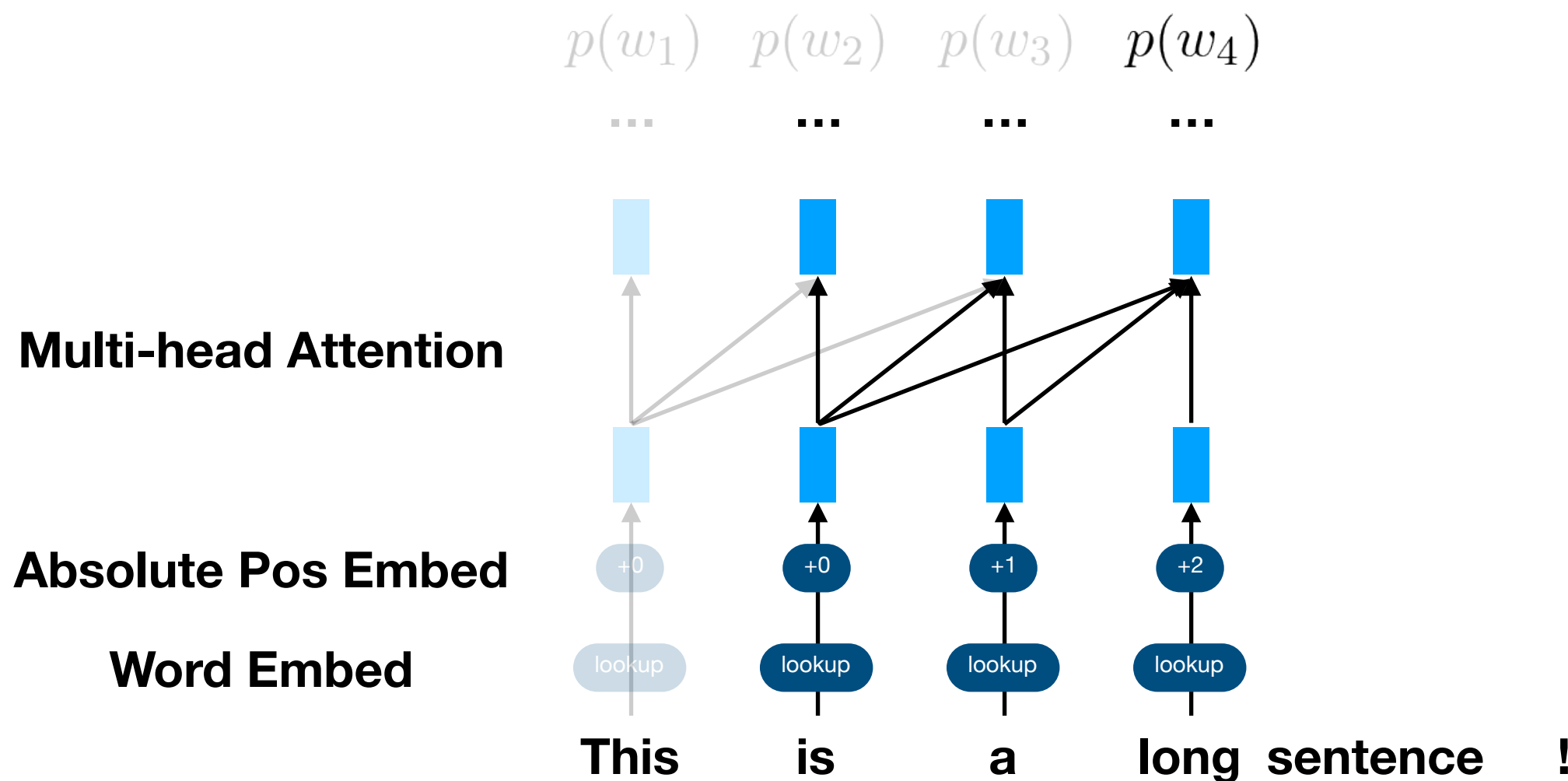
Transformer

- Transformer attention is parallel in time
- Absolute position embedding \rightarrow recomputation during inference



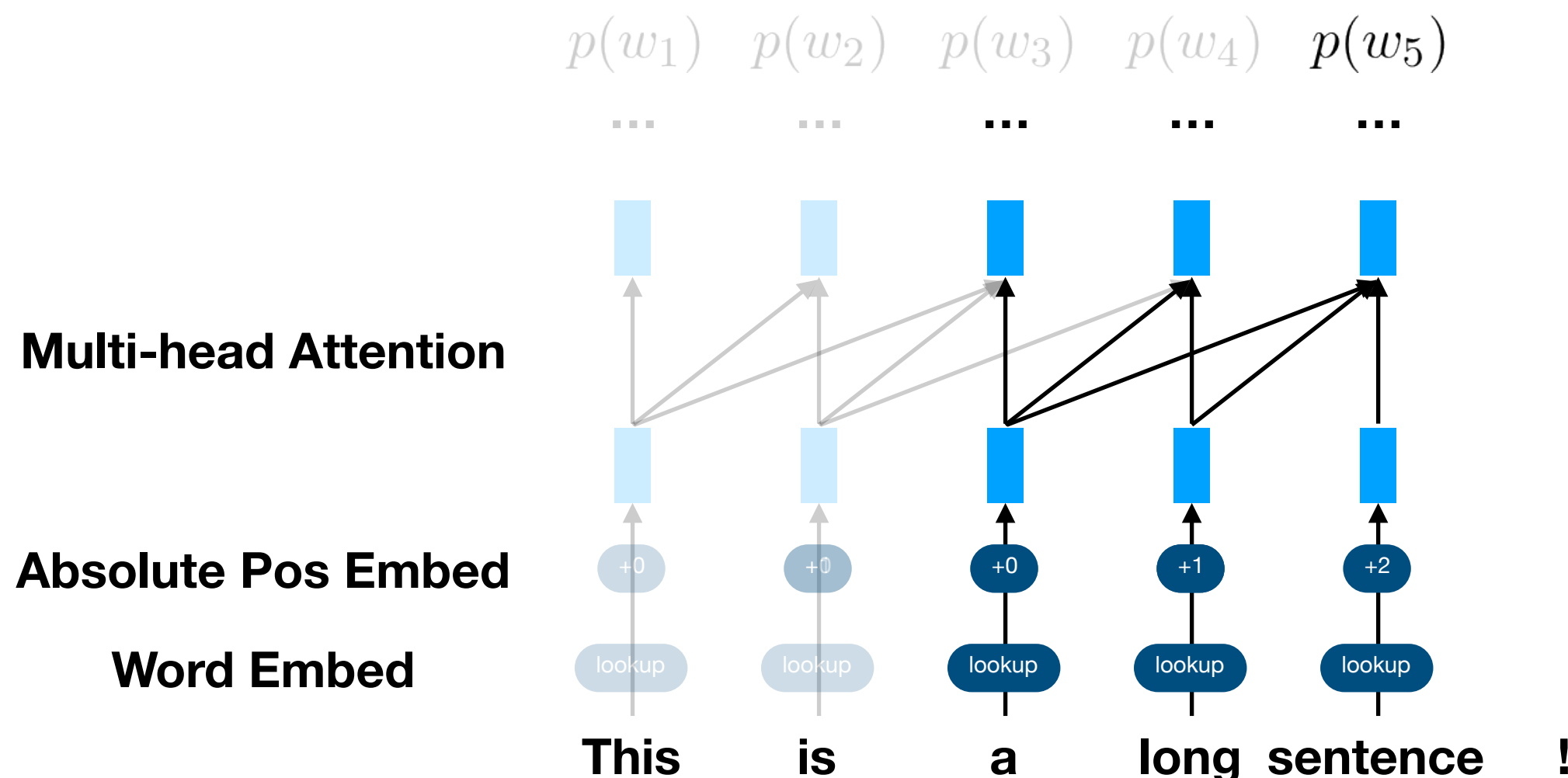
Transformer

- Transformer attention is parallel in time
- Absolute position embedding \rightarrow recomputation during inference



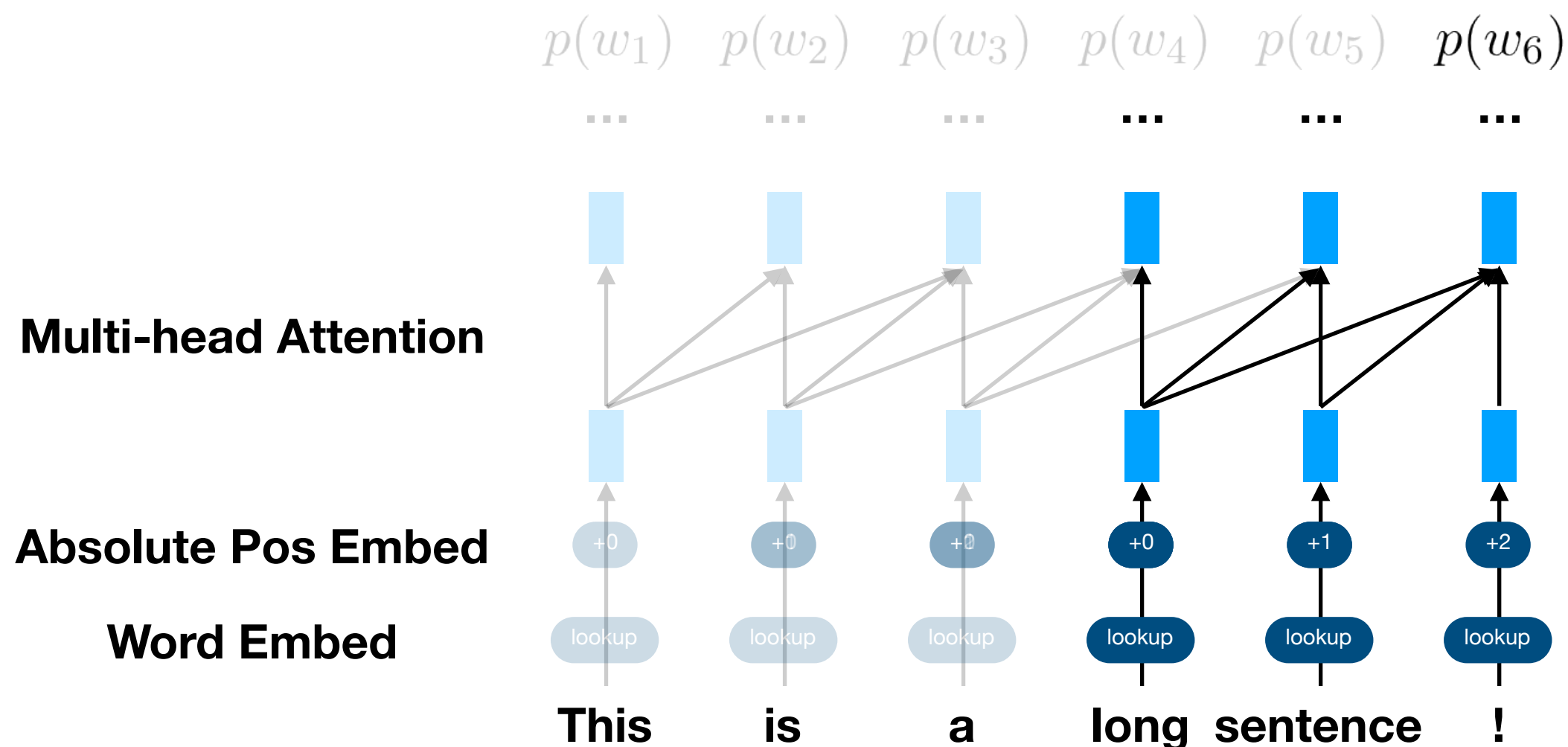
Transformer

- Transformer attention is parallel in time
- Absolute position embedding \rightarrow recomputation during inference



Transformer

- Transformer attention is parallel in time
- Absolute position embedding \rightarrow recomputation during inference



Transformer-XL

- Relative position embedding → efficient inference

This is a long sentence !

Transformer-XL

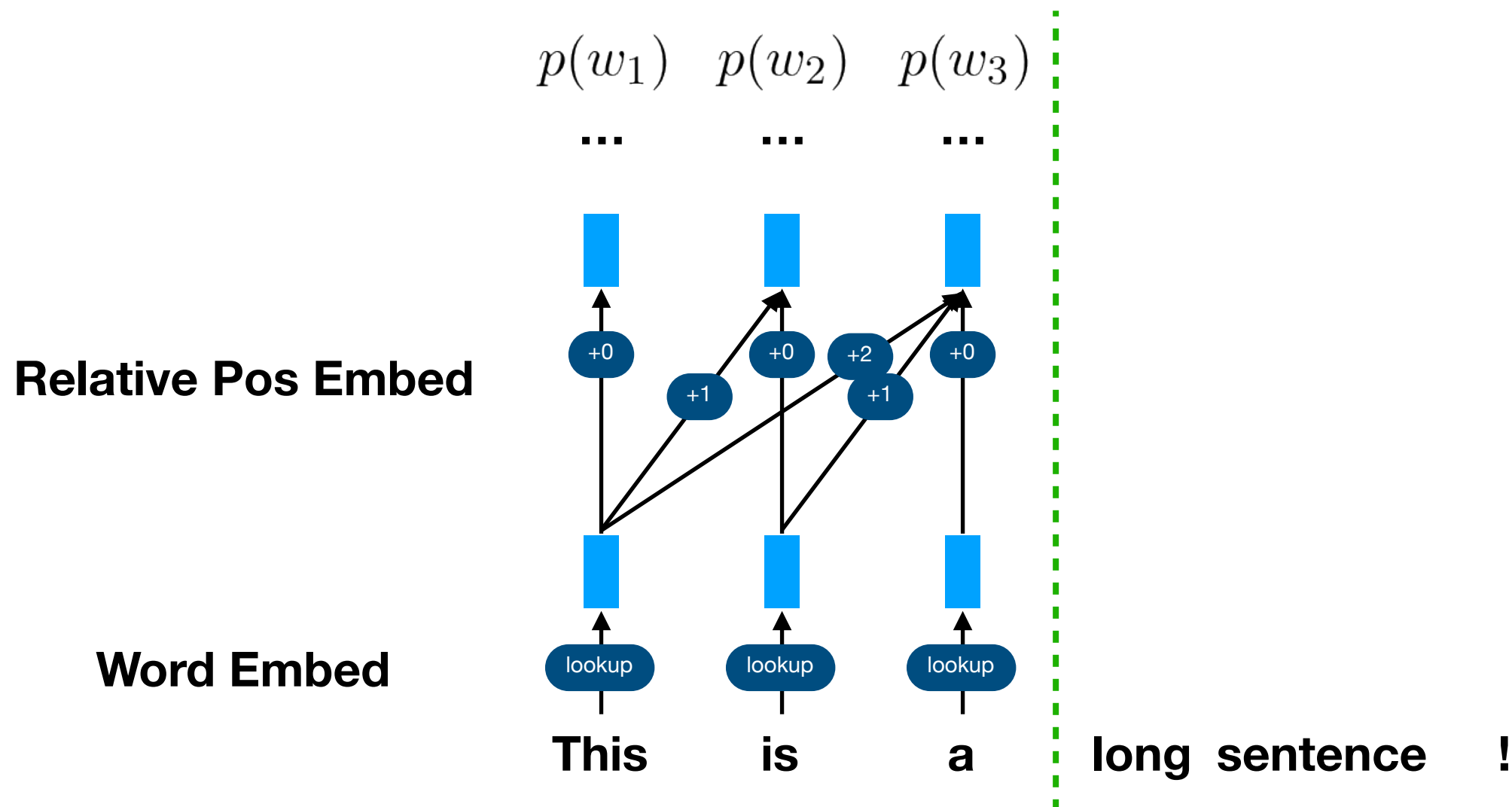
- Relative position embedding → efficient inference

This is a long sentence !



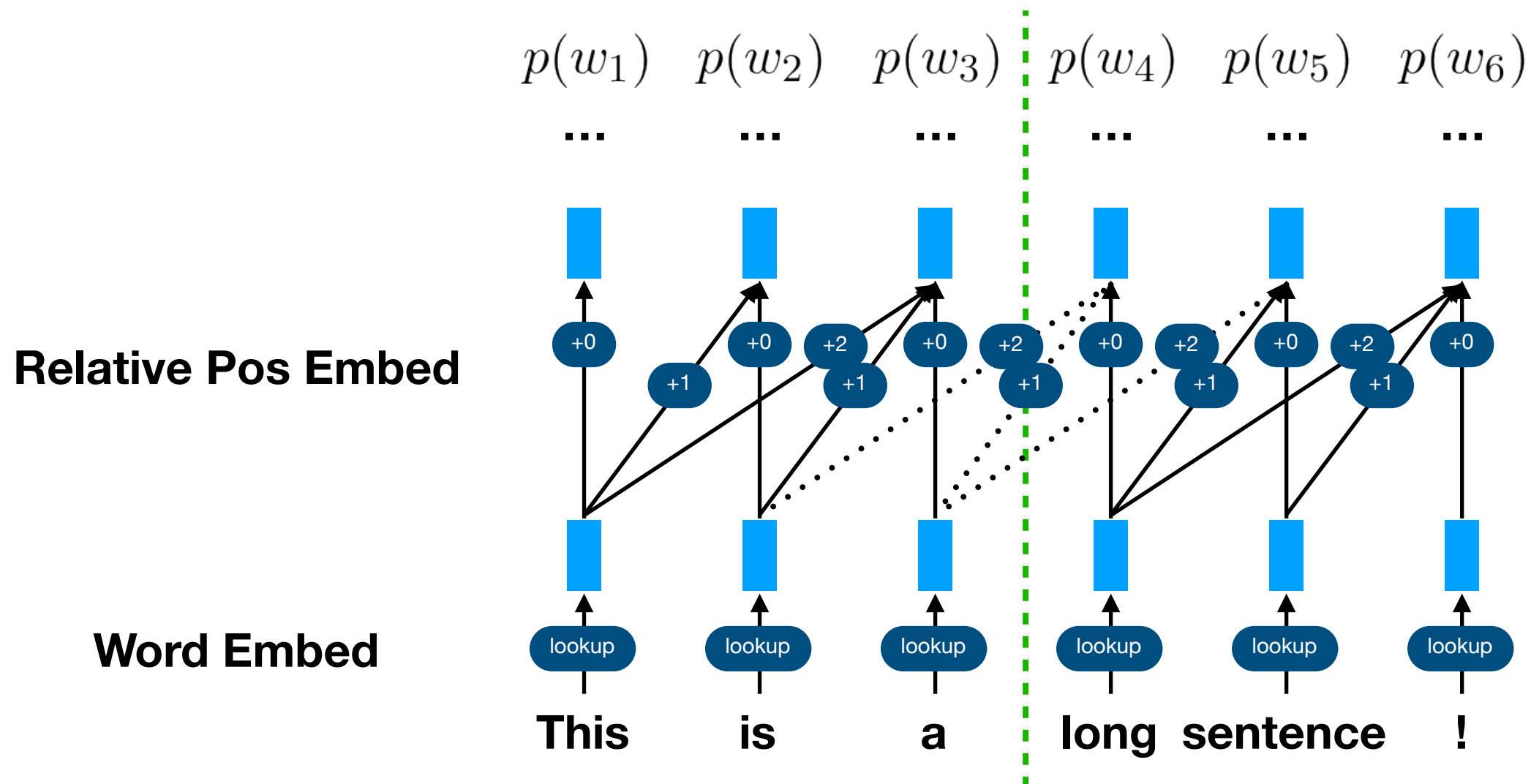
Transformer-XL

- Relative position embedding → efficient inference



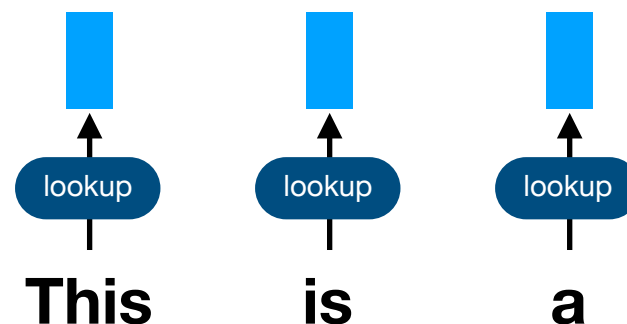
Transformer-XL

- Relative position embedding → efficient inference



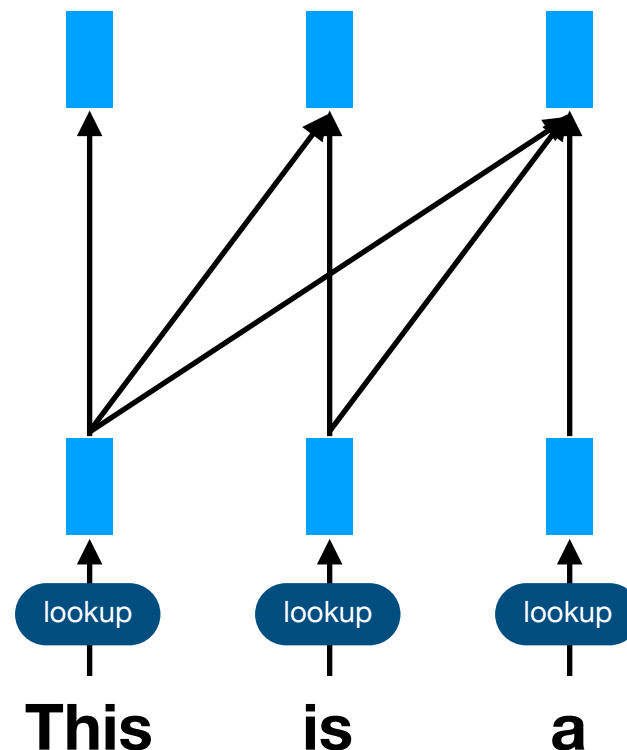
Weight-tied Nets

- E.g. Weight-tied transformer, Deep Equilibrium (DEQ) transformer
- Only one layer of weight
 - **But:** wider, more parameters and computation per layer



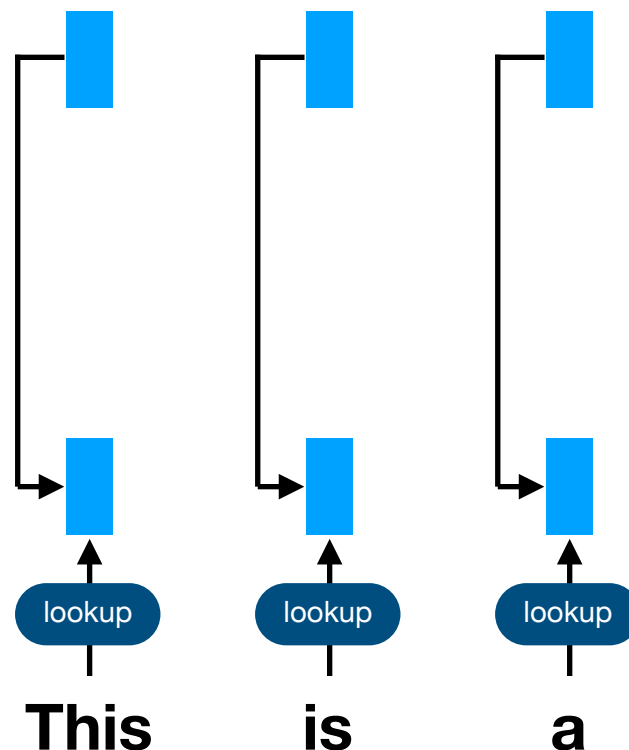
Weight-tied Nets

- E.g. Weight-tied transformer, Deep Equilibrium (DEQ) transformer
- Only one layer of weight
 - **But:** wider, more parameters and computation per layer



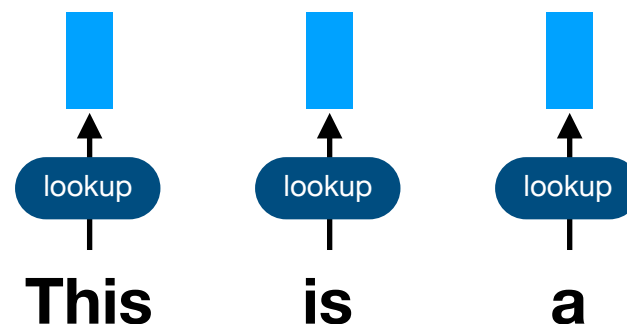
Weight-tied Nets

- E.g. Weight-tied transformer, Deep Equilibrium (DEQ) transformer
- Only one layer of weight
 - **But:** wider, more parameters and computation per layer



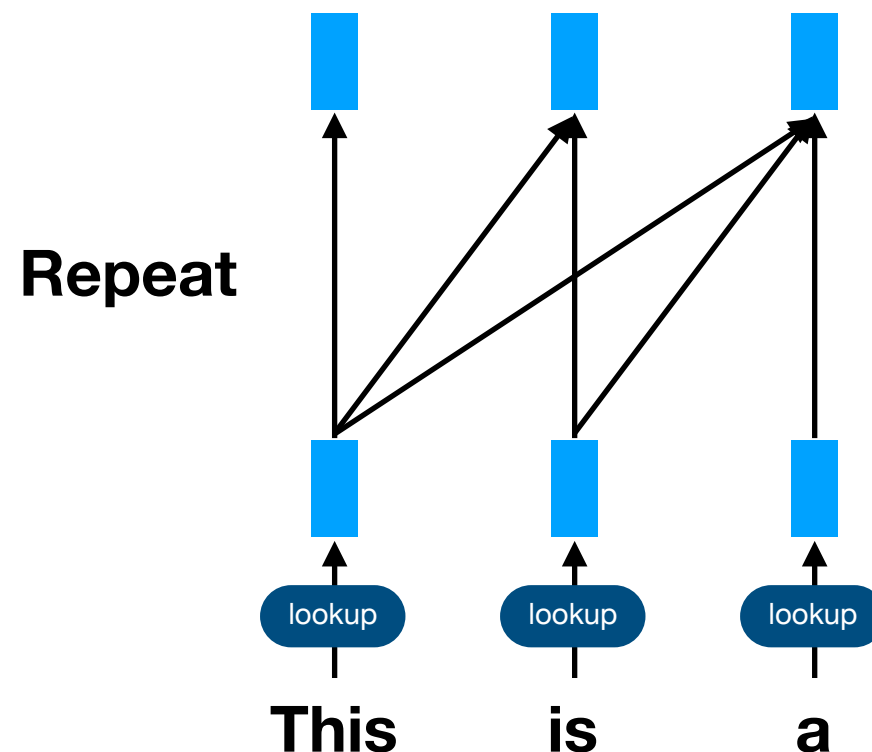
Weight-tied Nets

- E.g. Weight-tied transformer, Deep Equilibrium (DEQ) transformer
- Only one layer of weight
 - **But:** wider, more parameters and computation per layer



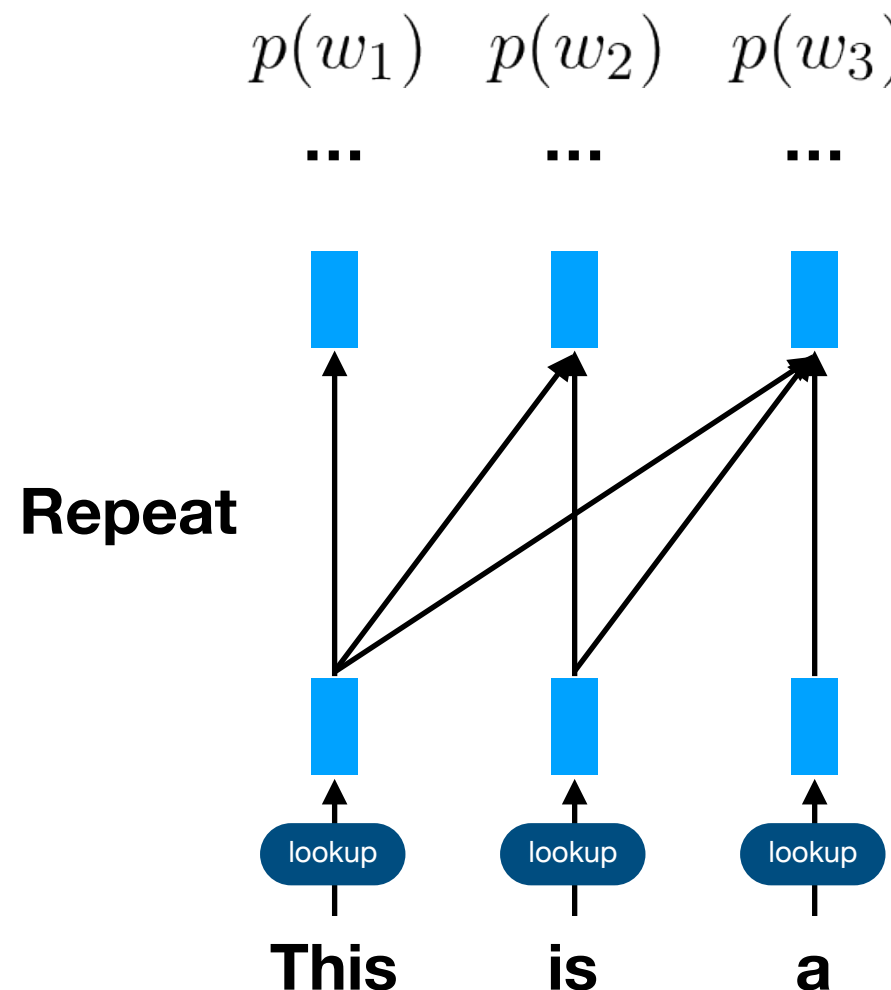
Weight-tied Nets

- E.g. Weight-tied transformer, Deep Equilibrium (DEQ) transformer
- Only one layer of weight
 - **But:** wider, more parameters and computation per layer



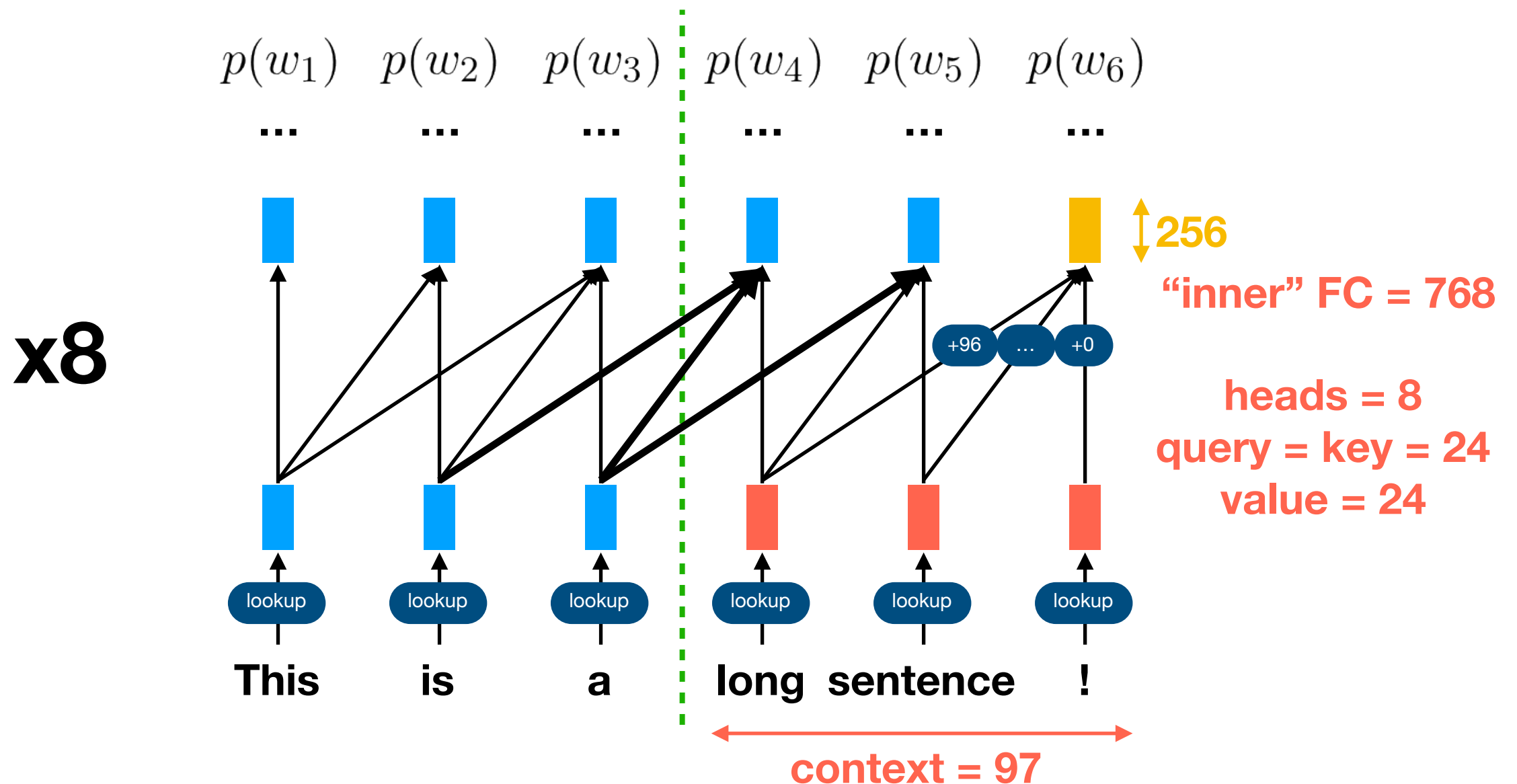
Weight-tied Nets

- E.g. Weight-tied transformer, Deep Equilibrium (DEQ) transformer
- Only one layer of weight
 - **But:** wider, more parameters and computation per layer



Our Base Model

- Fully differentiable transformer-XL

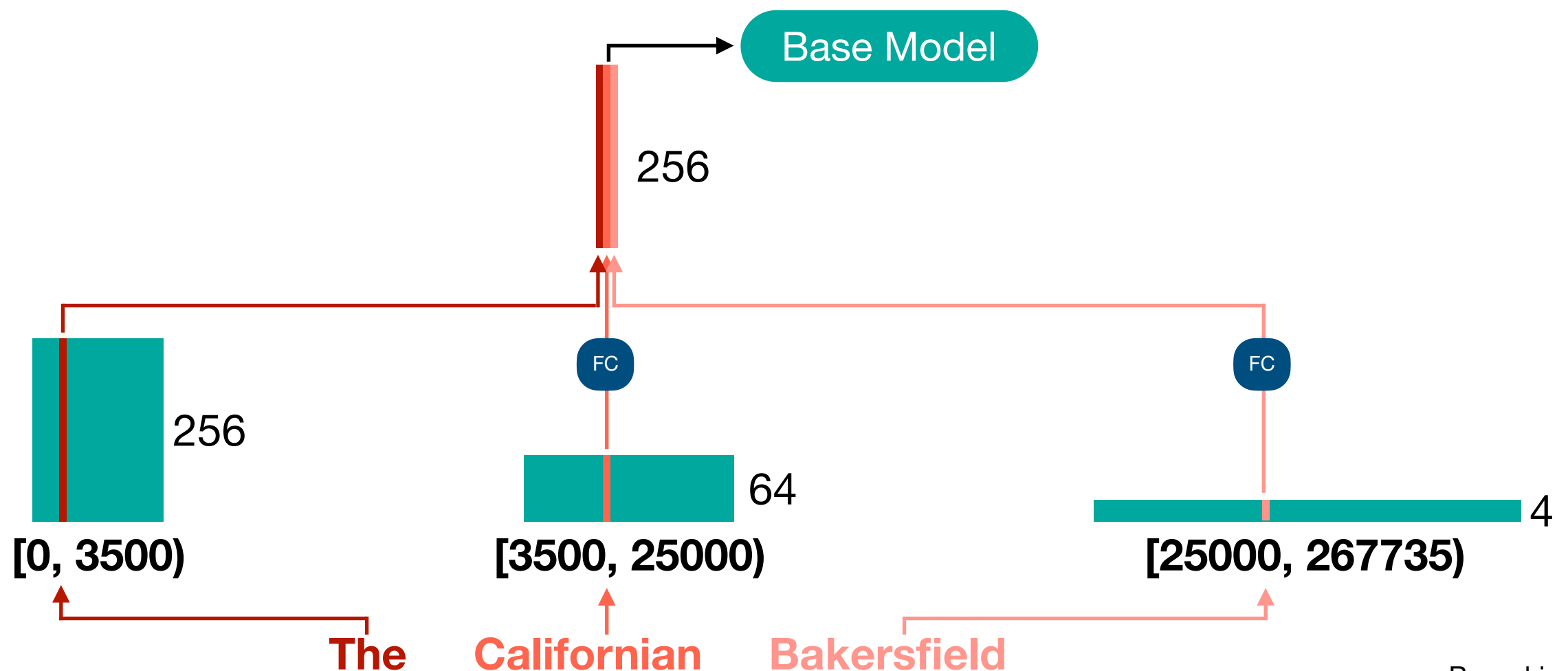


Observations

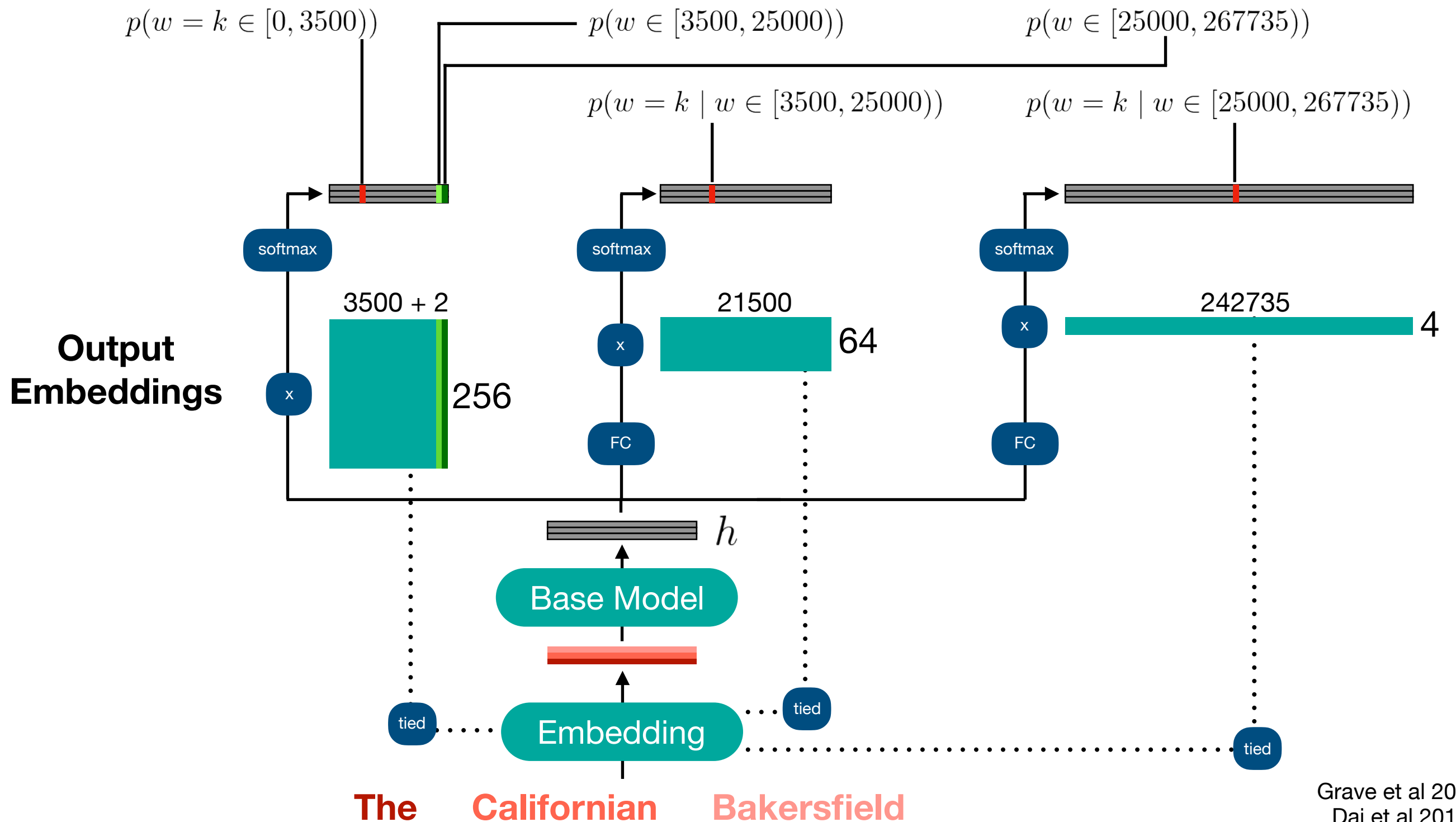
- Large vocab size: 267735
- Naive word embedding: $267735 \times 256 = 69\text{M}$ params
 - 138M mul/add's at output embedding
- Rare words need less representation!

Adaptive Word Embedding

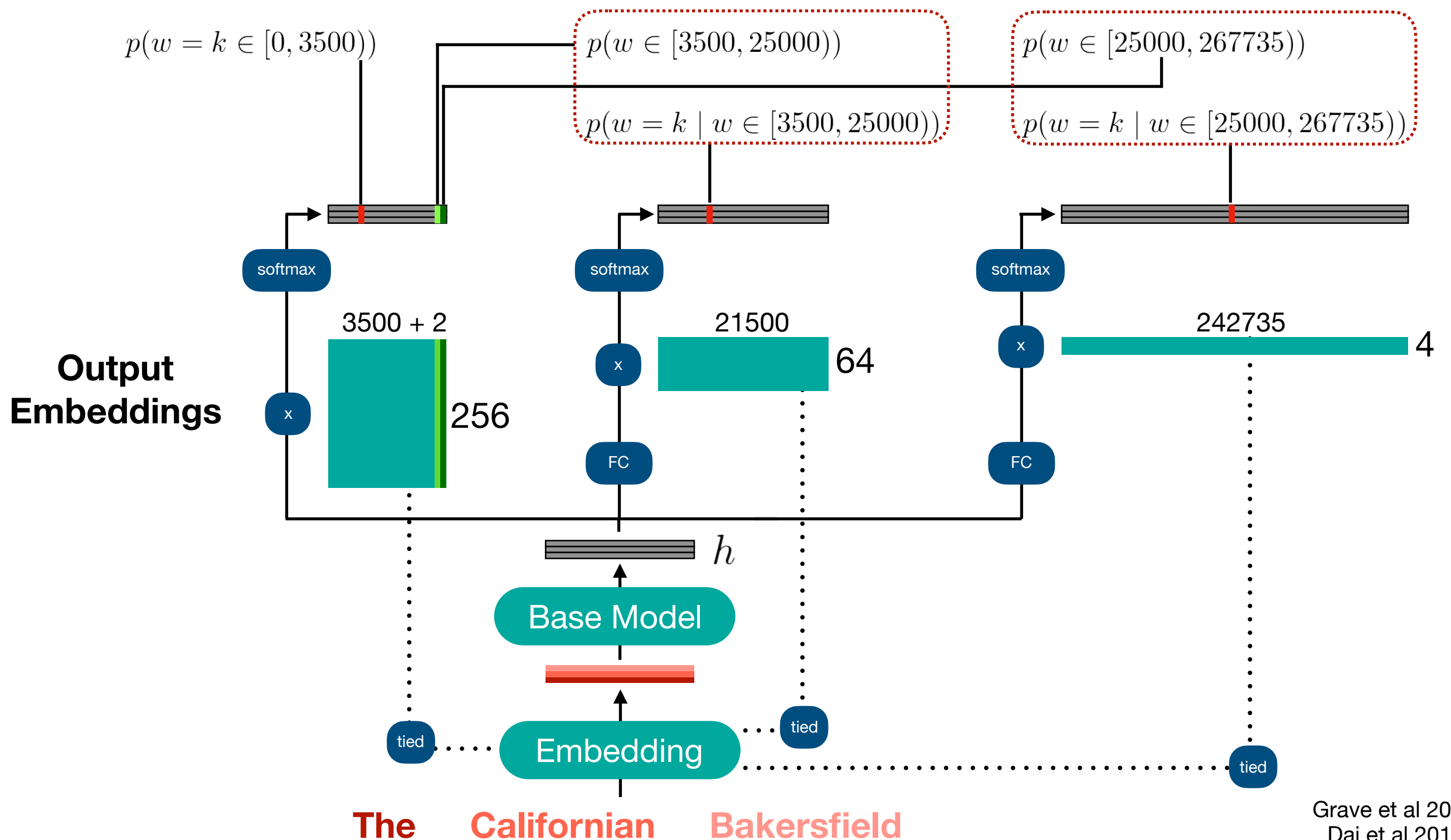
- Sort vocabulary and divide into 3 bins by frequency
- Low-rank decomposition



Adaptive Softmax



Adaptive Softmax



Performance

- Learning rate = $1e-3$
- Cosine learning rate decay
- 1 RTX 2080 Ti \rightarrow ~ 1 day
- Compute $\sim 2 * \text{Params}$

Layers	Embed Param (M)	Param (M)	Val PPL	Q, K, V	Inner
5	3.3	7.8	39.4	48	1024
6	3.3	8.3	37.9	32	1024
8	3.3	8.3	37.2	24	768
LSTM Baseline*	138	159	29.0		

* Baseline is Rae et al 2018

Observations

- Rare words have a much higher probability of recurring than occurring

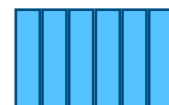
Vancouver was originally named **Gastown** and began as a settlement which grew around the site of a makeshift tavern on the western edges of Hastings Mill built on July 1, 1867, and owned by proprietor Gassy Jack. The original site is marked by the **Gastown** steam clock. **Gastown** then formally registered as a townsite dubbed Granville, Burrard Inlet.

- How to compensate for short attention context (97)?

Differentiable Cache

- Keep most recent activations in differentiable cache!
 - 0 parameters
 - Adds some computation

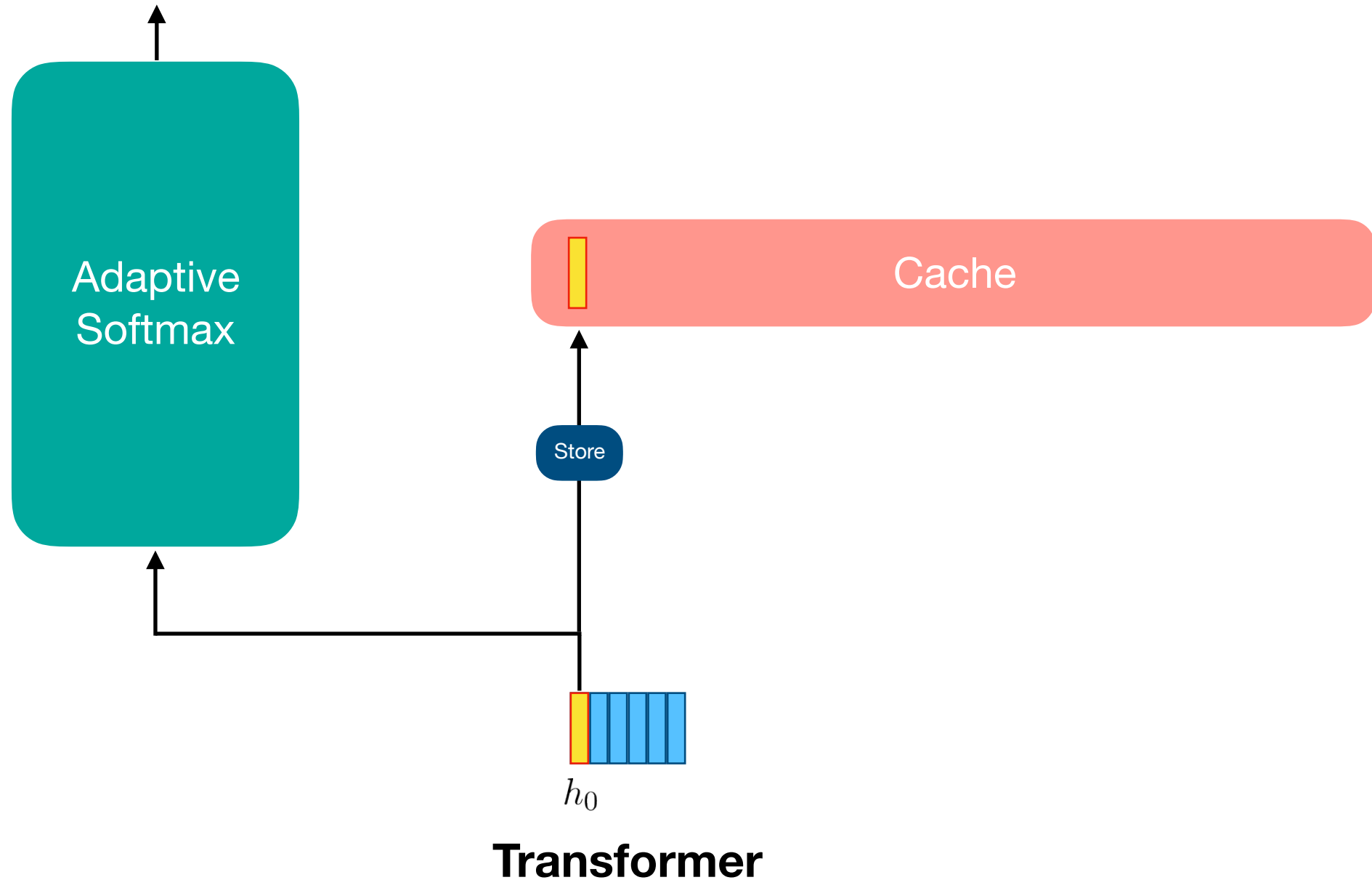
Differentiable Cache



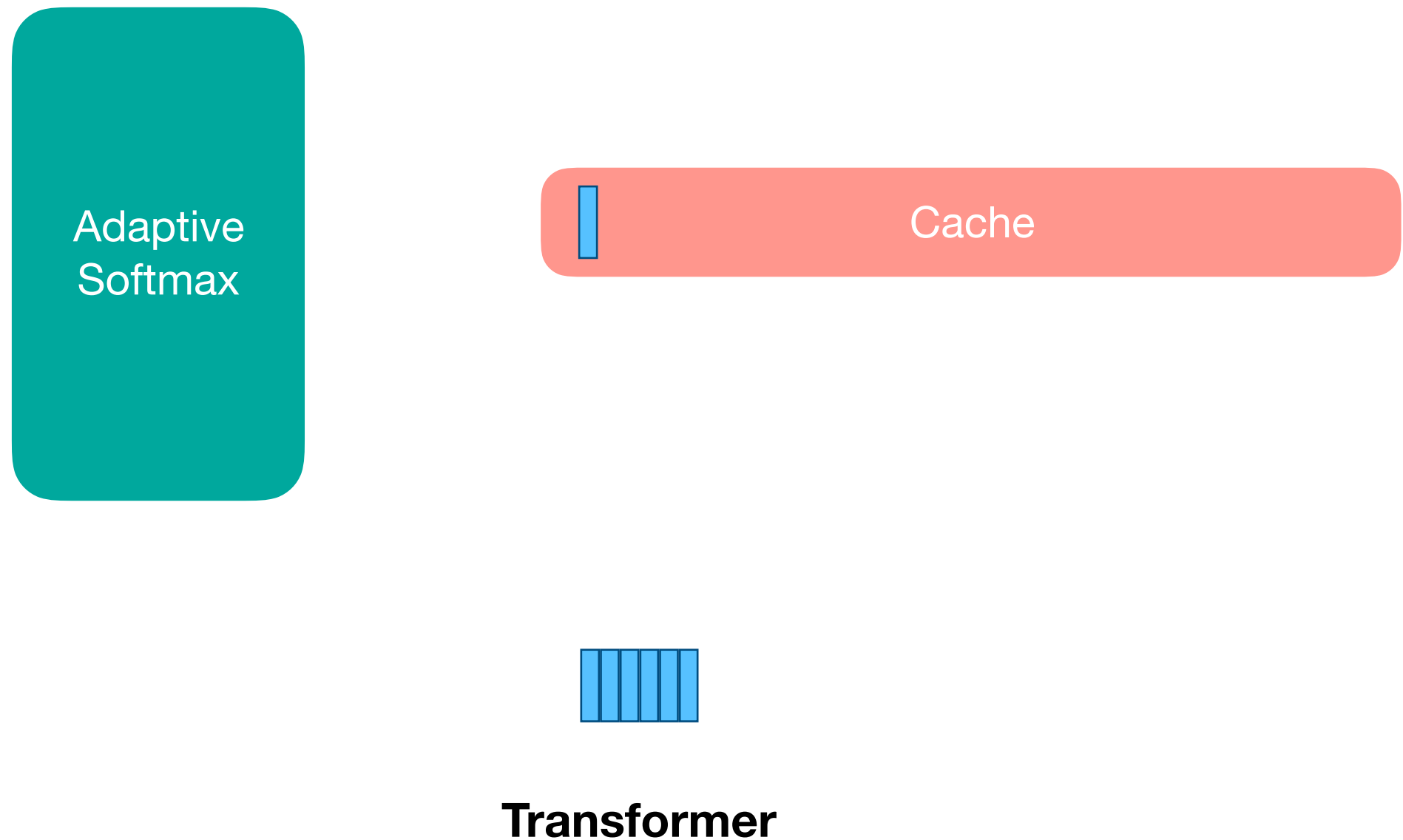
Transformer

Differentiable Cache

$$p(w_1) = p_{\text{embed}}(w_1)$$

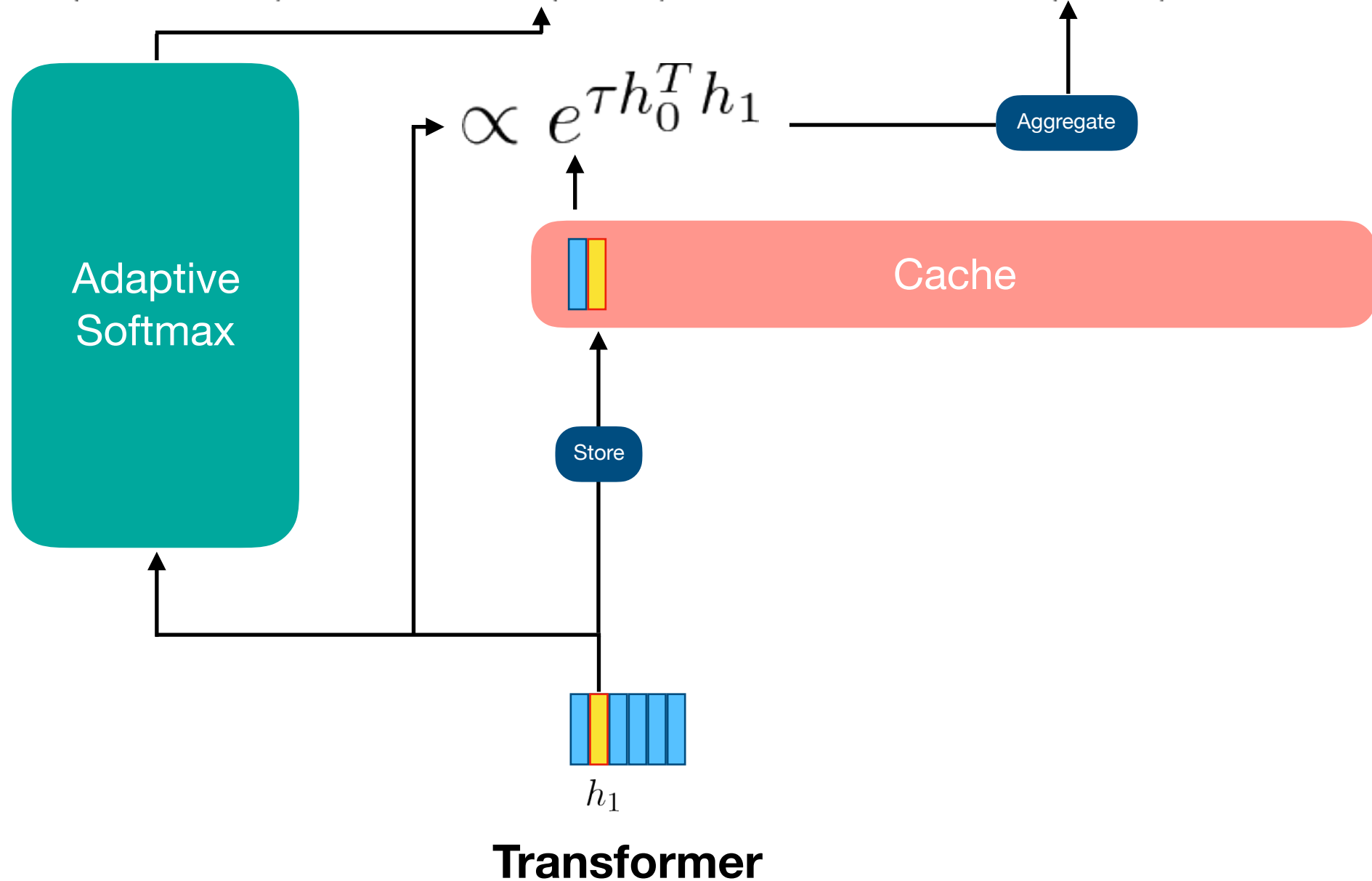


Differentiable Cache



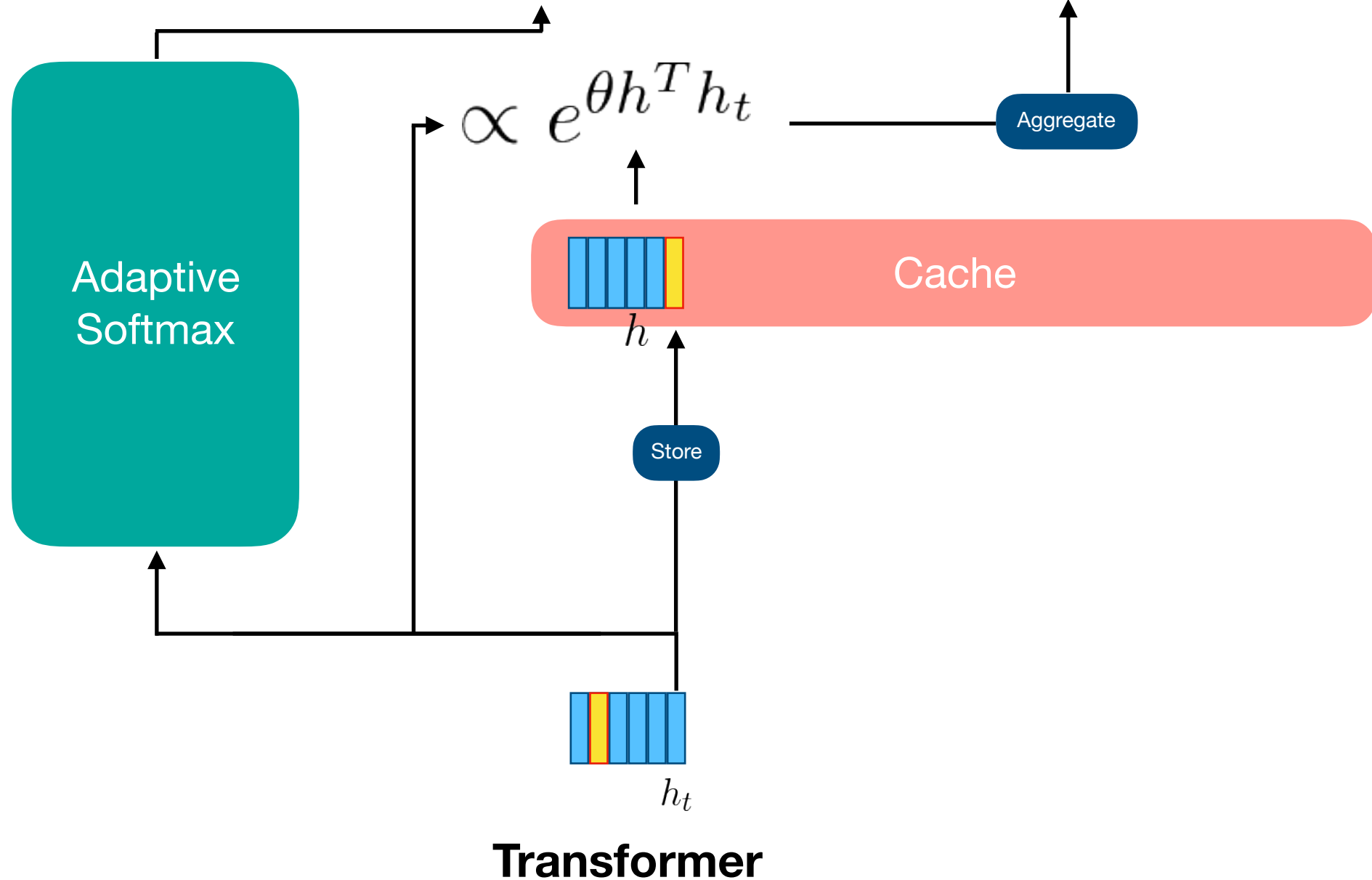
Differentiable Cache

$$p(w_2) = (1 - \lambda)p_{\text{embed}}(w_2) + \lambda p_{\text{cache}}(w_2)$$



Differentiable Cache

$$p(w_{t+1}) = (1 - \lambda)p_{\text{embed}}(w_{t+1}) + \lambda p_{\text{cache}}(w_{t+1})$$



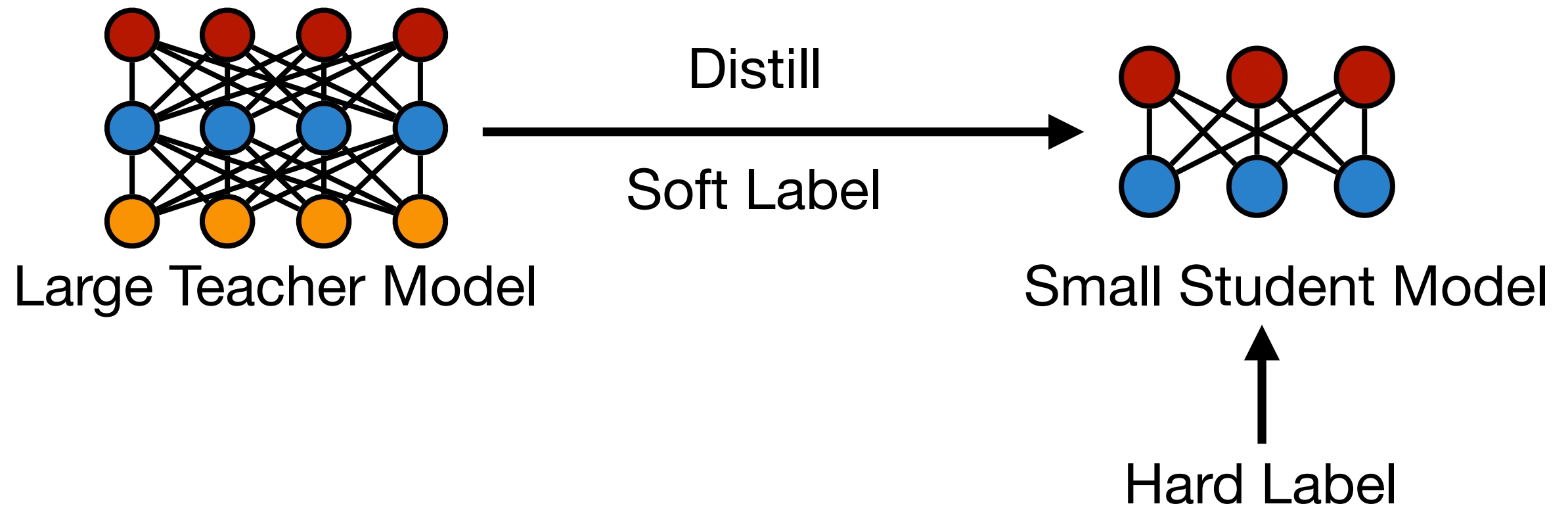
Performance

Param	Cache	Val PPL
8.3M	0	37.2
11.0M	0	35.3
15.2M	0	31.5
18.1M	0	30.1
8.3M	1000	33.1
8.3M	2000	32.0

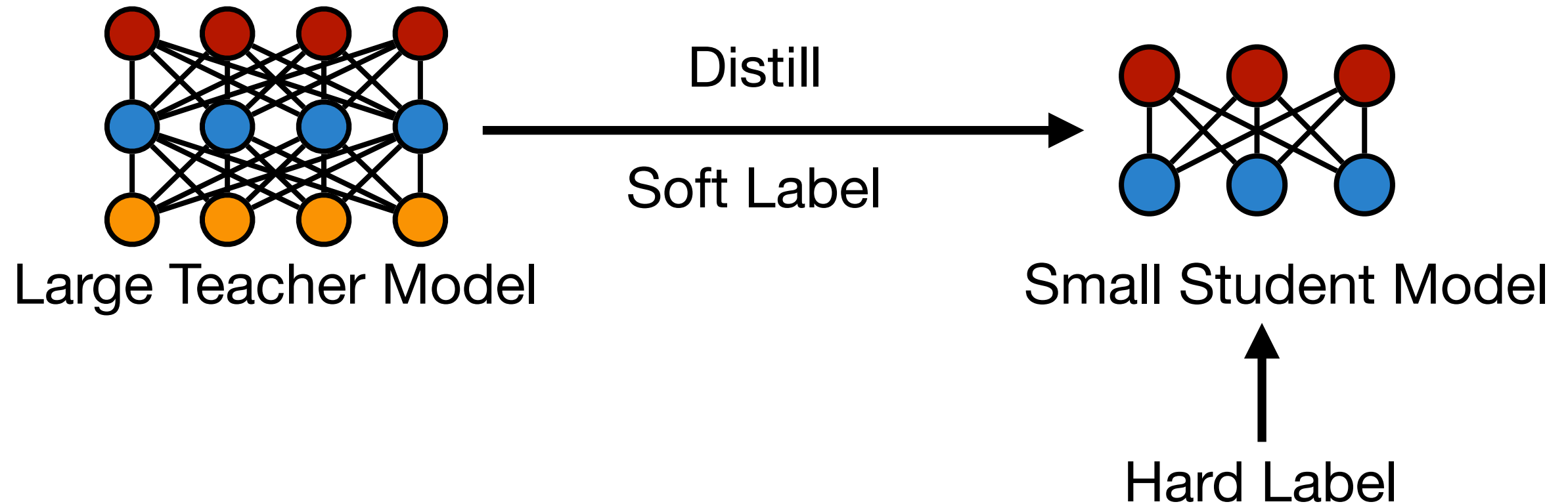
Model Compression

- Knowledge distillation
- Pruning
- Quantization

Knowledge Distillation

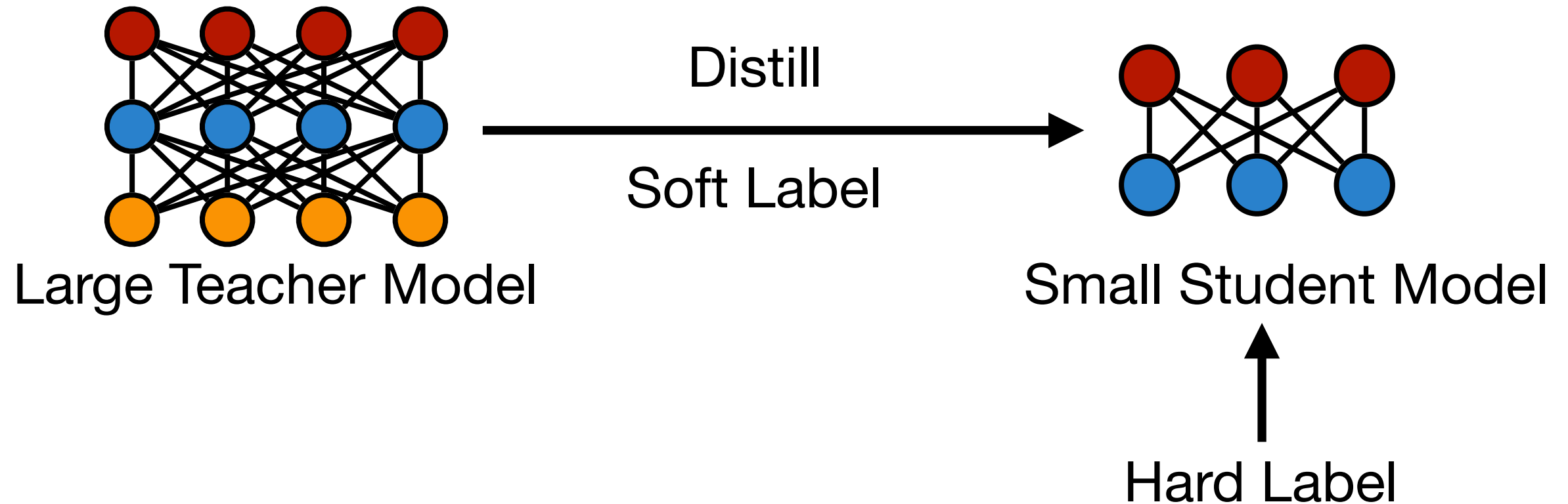


Knowledge Distillation



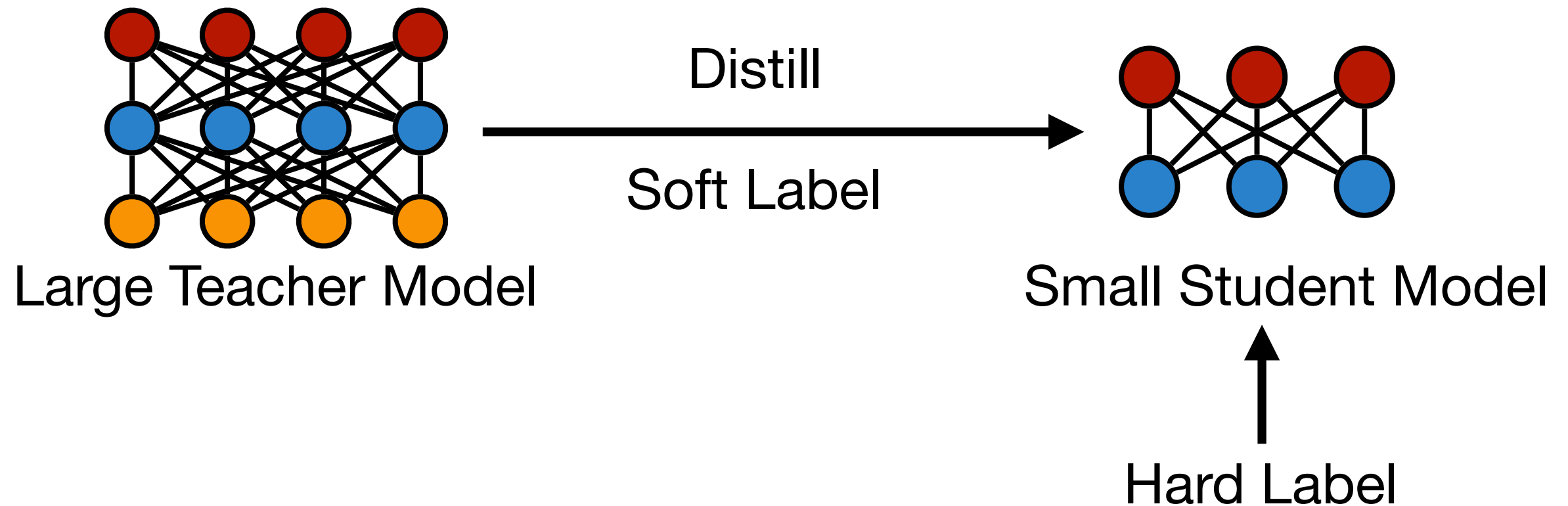
- Large teacher: 53M params, Test PPL 22

Knowledge Distillation



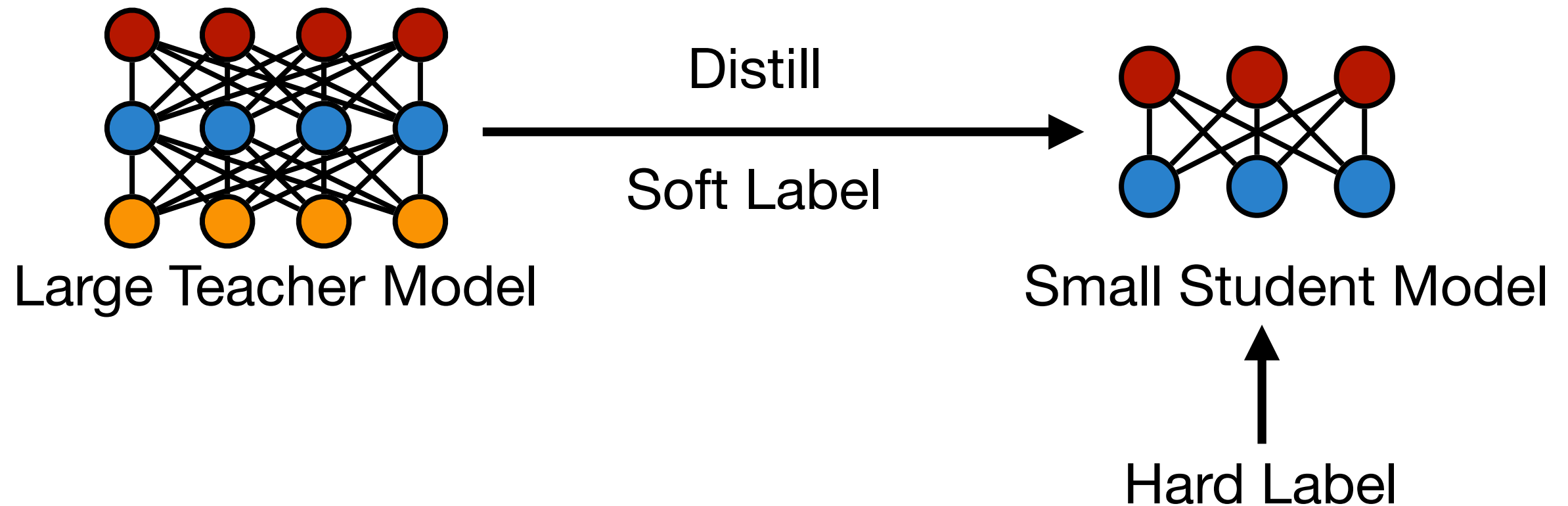
- Large teacher: 53M params, Test PPL 22
- Student: 8.3M params

Knowledge Distillation



- Large teacher: 53M params, Test PPL 22
- Student: 8.3M params
- Top 30 soft labels

Knowledge Distillation

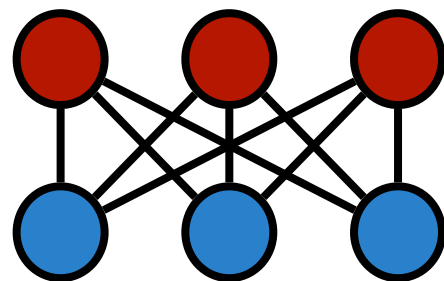


- Large teacher: 53M params, Test PPL 22
- Student: 8.3M params
- Top 30 soft labels
- Teacher annealing, first learn from teacher then learn from ground truth

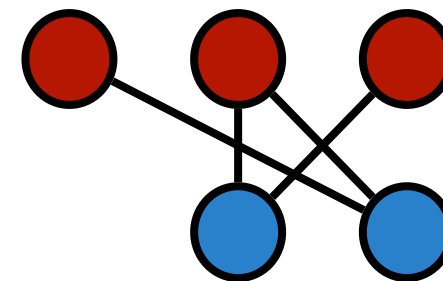
Performance

Model	Param (M)	Val PPL
Original	8.3	32.0
Distilled	8.3	30.7
Baseline	159	29.0

Pruning

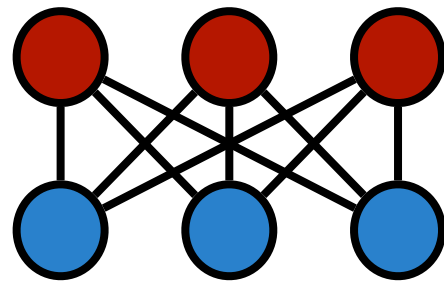


Trained Small Student Model

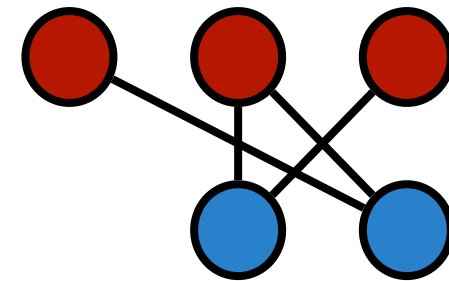


Pruned Small Student Model

Pruning



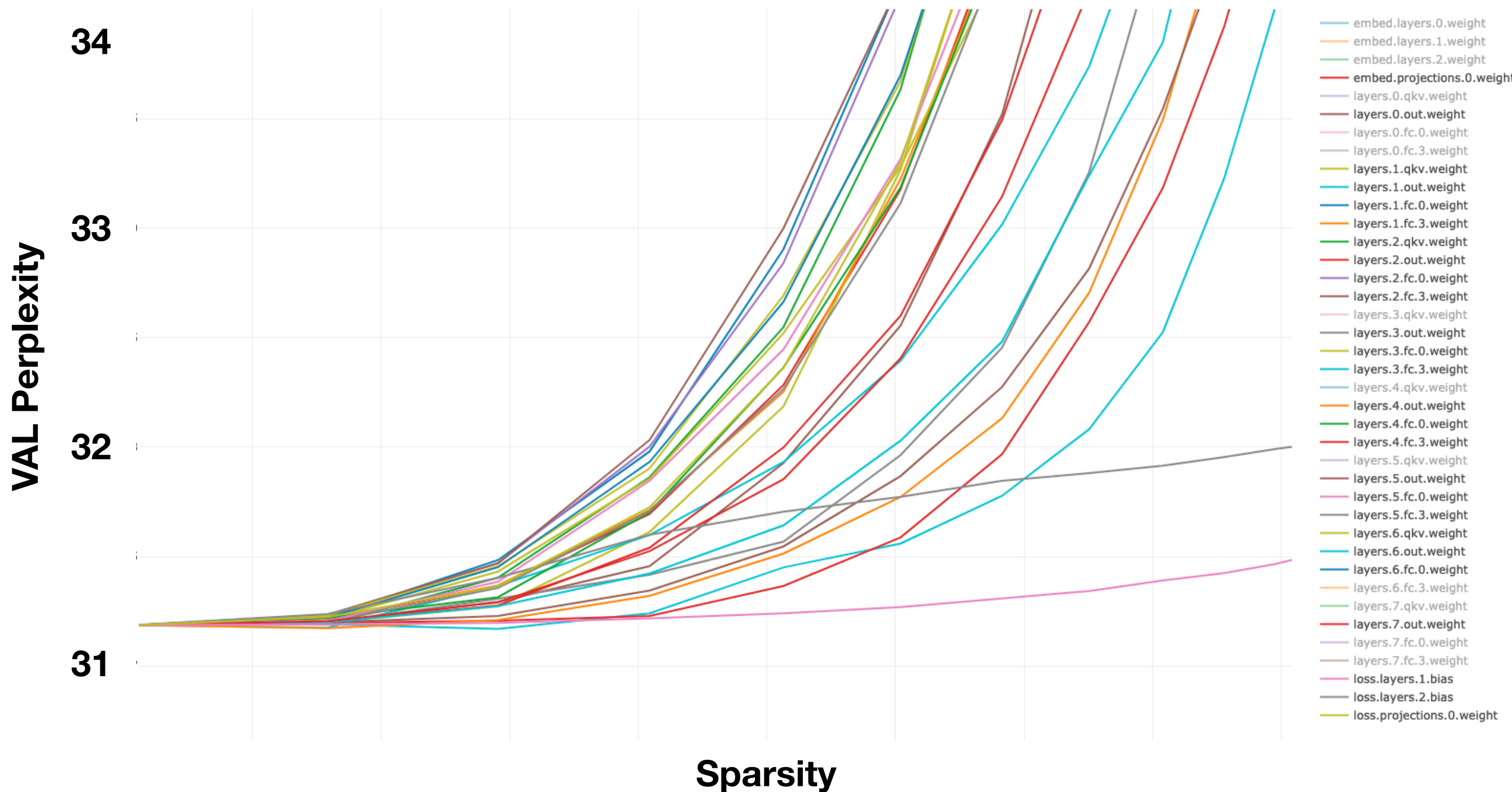
Trained Small Student Model



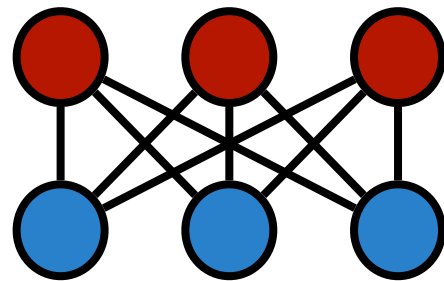
Pruned Small Student Model

- Sensitivity Check

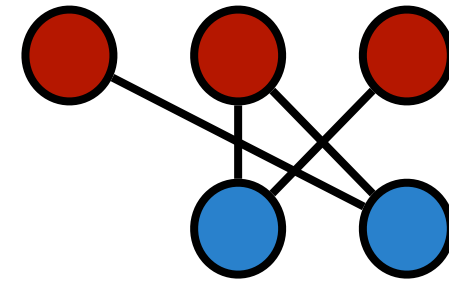
Sensitivity Check



Pruning



Pruning
→

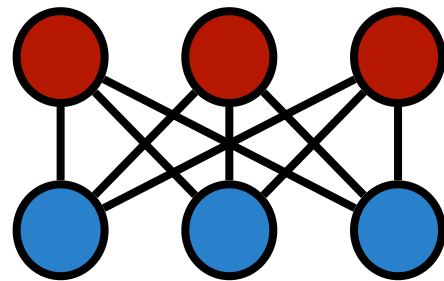


Trained Small Student Model

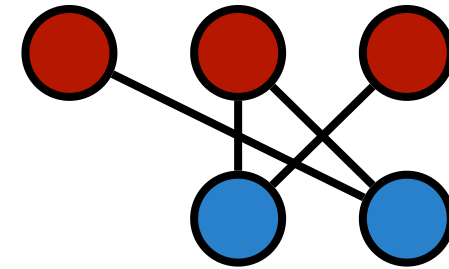
Pruned Small Student Model

- Sensitivity Check
- More sensitive layers will be pruned less

Pruning



Pruning
→



Trained Small Student Model

Pruned Small Student Model

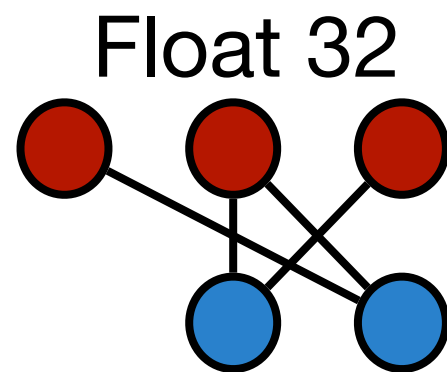
- Sensitivity Check
- More sensitive layers will be pruned less
- Automatic gradual pruning

Performance

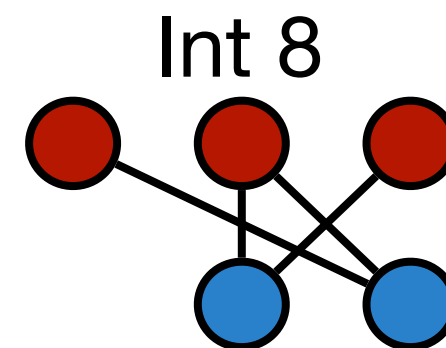
Model	Sparsity (%)	Param (M)	Compute (M)	Val PPL
1	42.1	5.06	11.8	34.3
2	40.1	5.22	12.1	34.0
3	33.9	5.74	13.1	34.3
Baseline	0	159	318	29.0

* Baseline is Rae et al 2018

Quantization

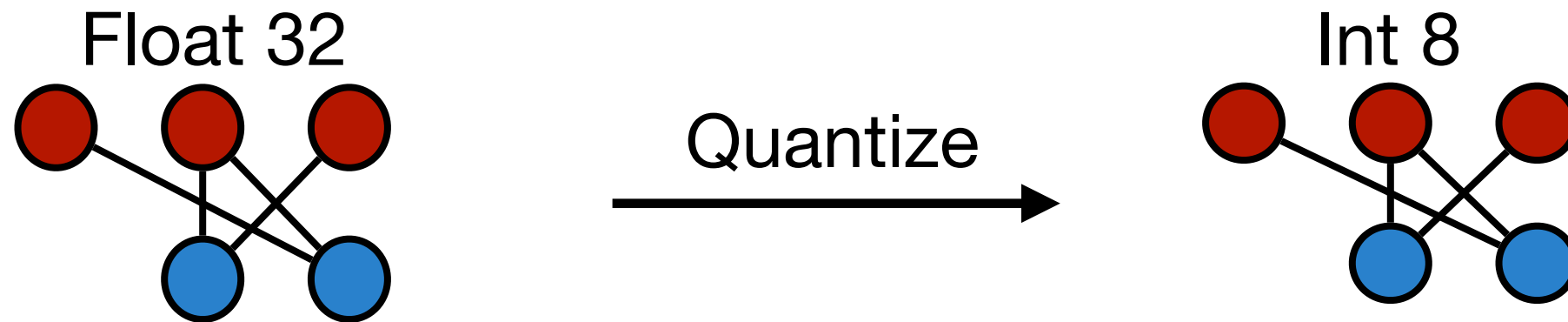


Pruned Student Model



Quantized & Pruned
Small Student Model

Quantization

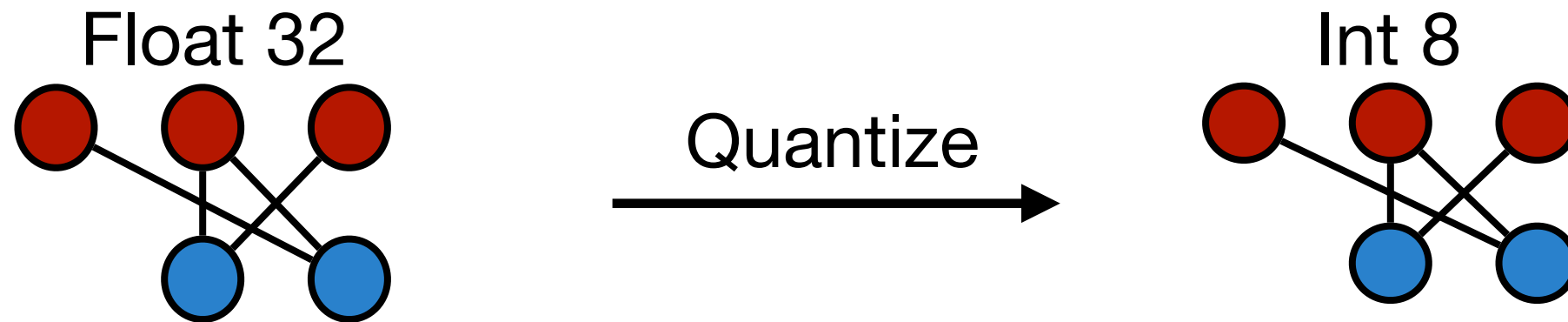


Pruned Student Model

Quantized & Pruned
Small Student Model

- Quantize weight, bias and activations of pruned model

Quantization

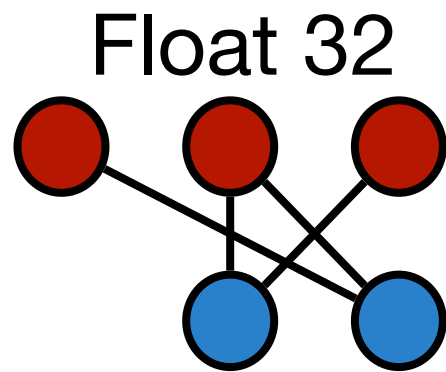


Pruned Student Model

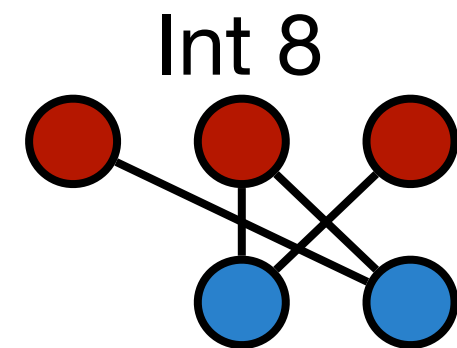
Quantized & Pruned
Small Student Model

- Quantize weight, bias and activations of pruned model
- Linear-range symmetric quantization

Quantization



Pruned Student Model



Quantized & Pruned
Small Student Model

- Quantize weight, bias and activations of pruned model
- Linear-range symmetric quantization
- Quantize to 8 or 9 bits based on trade-offs between quantization and pruning ratio

Results

Entry	Sparsity (%)	Quantize (bits)	32-bit Param (M)	32-bit Compute (M)	Val PPL	Test PPL	Score
1	42.1	9	1.61	7.83	34.3	34.95	0.0347
2	40.1	9	1.65	8.03	34.0	34.7	0.0356
3	33.9	8	1.63	8.48	34.3	34.95	0.0369
Baseline	0	16 (freebie)	~79.5	~239	29.0	29.2	~1.25

Score = param/159M + compute/318M

50x #params reduction and **31x** computation reduction than baseline

Open Source



GitHub open source of our models

Thank you!